# A Survey on Unstructured Document Annotation Using Content and Query Value Based

Ms. Pranjali Raut[1], Prof. Guraudev B. Sawarkar[2]

[1] M.Tech ( Department of Computer Science & Engineering )
VM Institute of Engineering and Technology,Nagpur

[2]Asst. Professor( Department of Computer Science & Engineering )
VM Institute of Engineering and Technology,Nagpur

**Abstract**— A large data is generated in different organization which is in textual format. This contain eloquent amount of structured information, which dwell underlying in the unstructured text. Document annotation is the task of adding metadata information in the document which is useful for information extraction. Many information mining algorithms facilitate the extraction of structured information from raw data but which is costly, inefficient and also shows contaminated results. We present, an adaptive technique that Collaborative Adaptive Data sharing platform (CADS) for document annotation and use of query workload to direct the annotation process. This paper proposes survey on facilitates the generation of structured metadata by identifying documents containing information of interest. Such information is further useful for querying the database. Our experimental assessment shows superior results compared to approaches that rely only on the textual content or only on the query workload, to recognize attributes of importance. In this paper, we are doing a survey on document annotation techniques.

**Keywords**— Annotation; CADS; METADATA; Adaptive technique; Structured & Unstructured.

## I. INTRODUCTION

Organizations generate huge amount of unstructured data. Advanced growth in data collection and storage technology made it possible to arrange this data at lower cost. Our goal is exploiting this stored data, in order to extract useful and actionable information. To get summarized search information is our requirement and to get this we arrange data in smart way. Annotation is one of the best techniques to arrange significant information and get effective search result. Annotations of documents are comments, notes, explanations, or other types of external remarks that can be attached to a web document or to a selected part of a document. As they are external, it is possible to annotate any web document independently, without needing to edit the document itself. From a technical point of view, annotations are usually seen as metadata, as they give additional information about an existing piece of data. Annotations of documents can be stored locally or in one or more database servers. When a document is searched, content of queries value each of these database servers, requesting the annotations related to that document in web server database. There are many annotation techniques are present that are based on attribute value pair. The strategies based on attribute value pair are effective method of document annotation. But there is restriction that document should be in structured format when using these systems. Also user has internal knowledge of attributes of document, as there are number of attributes because of them it will be difficult and infeasible to identify such attributes and its difficult approach to facilitate document

annotation. Hence algorithm should focus on those documents that contain words that are used during query. If we ignore contents of document then it will be unable to find out required information that's why document feature extraction is done on documents.

In this paper we propose cads (collaborative adaptive data sharing platform) which are an "annotation as user created". It assists fielded data annotations. In this the direct use of key contribution for the query workload by using direct annotation process and also examining the content of the document. Our aim is to prioritize the annotation of documents towards generating attribute names and attribute values for attributes that will often used by querying users and these attribute values will provide best possible results to the user where in users will have to deal only with relevant results.

## II. RELATED WORK

S.R. Jeffery, M.J. Franklin, and A.Y. Halevy [1]: proposed a paper Pay-as-You-Go User Feedback for Dataspace Systems. This system propose a system which is a line of work towards using more expressive queries that leverage annotations is the "pay-as – you – go " querying strategy in data -spaces. In data spaces users provide data integration hints at querying time. But in this paper it is assumed that data sources already contain structured information and the problem is to match the query attributes with the source attribute.

K. Saleem, S. Luis, Y. Deng, S.-C. Chen, V. Hristidis, and T. Li [2]: proposed a paper "Towards a Business Continuity Information Network for Rapid Disaster Recovery In this paper they consider the Crisis Management and Disaster Recovery have gained immense importance in the wake of recent man and nature inflicted calamities. They proposed a solution or model for pre-disaster preparation and post disaster business continuity/rapid recovery. In case of disaster need of rapid information retrieval and sharing increases. This paper proposed a disaster management model which works good at some extent but it is not considering the effective retrieval.

R.T. Clemen and R.L.Winkler[3]: proposed a paper "Unanimity and Compromise among Probability Forecasters" In this paper they work on probabilities of particular uncertain event. This helps us to find out annotation and attributes.

M. Franklin, A. Halevy, and D.Maier[4]: proposed a paper "From Databases to Dataspaces: A New Abstraction for Information Management ".It proposed a solution to Laplace smoothing to avoid zero probabilities for the attributes that do not appear in the workload. It helps us to converge towards accuracy.

G. Tsoumakas and I. Vlahavas [5]: propose a paper Random K-Labelsets: An Ensemble Method for Multilabel Classification. This paper proposes an ensemble method for multilabel classification. The RAndom k-labELsets (RAKEL) algorithm constructs each member of the ensemble by considering a small random subset of labels and learning a single-label classifier for the prediction of each element in the power set of this subset. In this way, the proposed algorithm aims to take into account label correlations using single-label classifiers that are applied on subtasks with manageable number of labels and adequate number of examples per label. Using this we can take into account the correlation between tags for annotations. But in this collaborative annotation is missing.

P. Heymann, D. Ramage, and H. Garcia-Molina [6]: proposed a paper "Social Tag Prediction": This paper give solution for prediction of tags for particular object. We can adopt this for out suggesting annotation concept.

### III. PROPOSED SCHEME

This system suggest, Collaborative Adaptive Data Sharing platform (CADS). CADS are annotate-as-you-create infrastructure that facilitates fielded data annotations. The goal of CADS is to lower the cost creating annotated documents that can be useful for commonly issued semi structured queries. [Figure-1] represents work flow of CADS. The CADS system has two types of actors: producers and consumers. Producers upload data in the CADS system using interactive insertion forms and consumers search for relevant information using adaptive query forms.
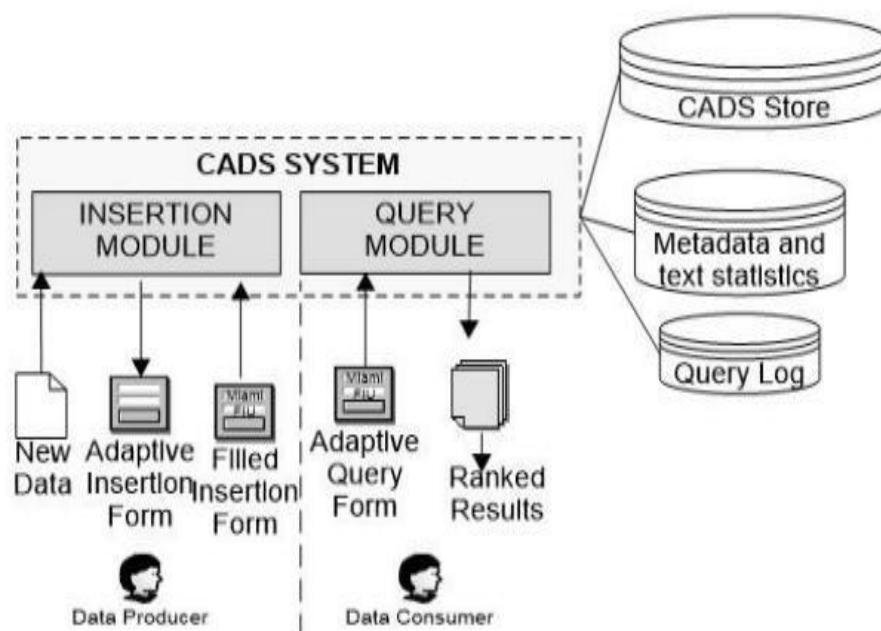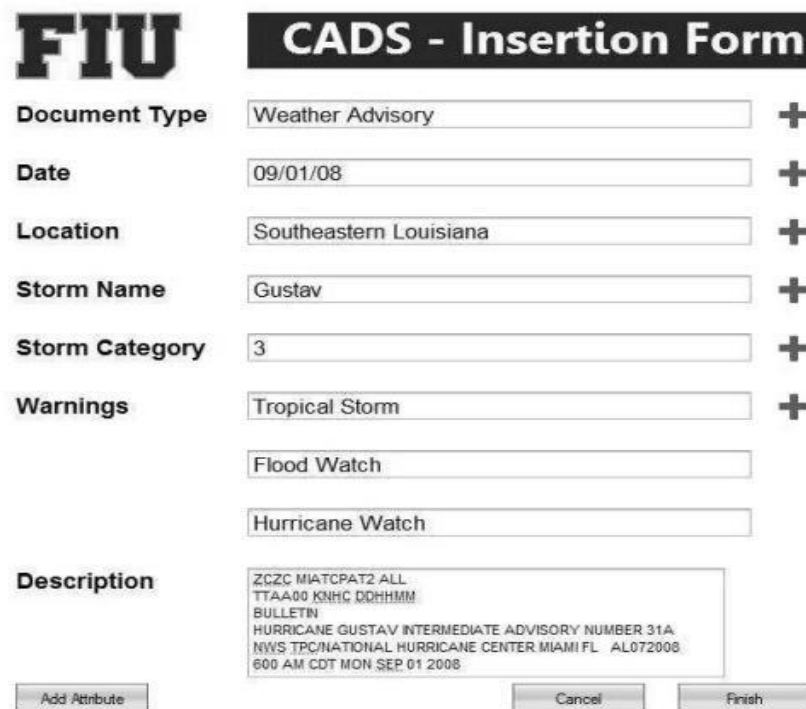
Fig: 1. CADS Workflow

In proposed system, the author generates a new document and uploads it in storage. After uploading the document, CADS analyses the text and creates adaptive insertion form as shown in **[Figure-2]**. The form contains the best attribute names which are present in the document and information needed for query workload and most probable values of the attributes given in the document. The author has ability to check the form, modify the metadata if it is necessary and finally submit the document for storage.



Fig: 2. Adaptive insertion form

While retrieving attribute names, the adaptive insertion form also retrieve the attribute values by using IE (Information Extraction) Algorithm. In order to retrieve contains of the text file information extraction (IE) algorithm is used.

### 1. *Information Extraction Algorithm:*

Step 1: Select a text file for extraction.
Step 2: Parse the text file. Ignore stop words from it and count frequency of high querying keywords which will be important for content based search. Maintain frequency count of

these keywords appearing in only single document.

Step 3: Upload the file on server.

Step 4: Then fill all the annotations which are relevant to the document which can be useful for query based searching.
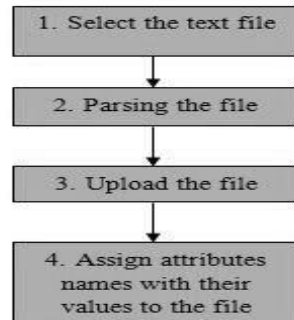


Fig: 3. IE Algorithm

The key contribution of this work is the "attribute suggestion" problem, which accounts for the query workload, and identifies the attributes that are present in the document, but not their values. There are two exclusive properties for indentifying and suggesting attributes for a document d.

- The attribute must have high querying value (QV) with respect to the query workload W.
- The attribute must have high content value (CV) With respect to d.

## 2. QV, CV Computation and Combining Algorithm:

Step 1: Enter the queries for retrieving the document Example: location='Mumbai' and year=2012.

Step 2: Disjoin the queries and pass it to database for retrieving.

Step 3: Check all annotated results and show the related results to user.

Step 4: For much efficient and accurate results, users should try to enter maximum queries they can.

## 3. Modules:

1. Registration
2. Login
3. Document Upload
4. Search Techniques
5. Download Document

## IV. CONCLUSION

This paper surveys work related to document annotation using content and querying value. This project is proposing adaptive methods to suggest relevant, recommended attributes to

annotate a document while also trying to satisfy the user querying, searching needs. Our solution is based on a probabilistic framework that views the confirmation in the document content and the query workload. The main advantage of our application is mainly that when users perform query based search, they could get minimum and distinct outcomes/results where it could be easy for retrieval. Experiments shows using these techniques, workload of application can reduce by large amount. It improves the annotation process and visibility of documents.

## REFERENCES

[1]     Eduardo J.Ruiz, Vagelis Hristidis , Panagiotics G.Iperiotis, "Facilitating Document Annotation using Content and            Querying Value" IEEE Transaction on Knowledge and data engineering 2014.

[2]     S.R. Jeffery, M.J. Franklin, and A.Y. Halevy, "Pay-as-You-Go User Feedback for Data space Systems," Proc.ACM SIGMOD  Int'l Conf. Management Data,2008.

[3]     K. Saleem, S. Luis, Y. Deng, S.-C. Chen, V. Hristidis, and T. Li, "Towards a Business Continuity Information Network for Rapid Disaster Recovery," Proc. Int'l Conf. Digital Govt. Research (dg.o '08), 27.

[4]     A. Jain and P.G. Ipeirotis, "A Quality-Aware Optimizer for Information Extraction," ACM Trans. Data base Systems, vol. 34, article 5, 2009.

[5]     J.M. Ponte and W.B. Croft, "A Language Modelling Approach to the Information Retrieval," Proc. 21st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '98), pp. 275-281,  http://doi.acm.org/10.1145/290941.291008

[6]     J. Madhavan et al., "Web-Scale Data Integration: You Can Only Afford to Pay as You Go," Proc. Third Biennial Conf. Innovative Data Systems Research (CIDR),  2007.

[7]     A. Halevy, Z. Ives, D. Suciu, and I. Tatarinov, "Schema Mediation in Peer Data Management Systems," Proc. 19th Int'l Conf. Data Eng., pp. 505-516,  Mar. 2003.

[8]     Microsoft, Microsoft Sharepoint, http://www.microsoft.com/sharepoint/, 2012.

[9]     SAP, Sap Content Manager, https://www.sdn.sap.com/irj/sdn/nw-cm, 2011.

[10]    O. Etzioni, M. Banko, S. Soderland, and D.S. Weld, "Open Information Extraction from the Web," Comm.  ACM, vol. 51,pp. 68-74,  http://doi.acm.org/10.1145/1409360.

[11]    C. D. Manning, P. Raghavan, and H. Schütze, Introduction to Information Retrieval, 1st ed. Cambridge University Press,July 2008.

[12]    R. Fagin, A. Lotem, and M. Naor, "Optimal aggregation algorithms for middleware," J. Comput. Syst. Sci., vol. 66, pp. 614–656,  June 2003.

[13]    M. J. Cafarella, J. Madhavan, and A. Halevy, "Web-scale extraction of structured data," SIGMOD Rec., vol. 37, pp. 55–61,  March 2009.

[14]    "Google, "Google Base, http://www.google.com/base, 2011..